

## Introduction

- Syntax: the form or structure of the expressions, statements, and program units
- Semantics: the meaning of the expressions, statements, and program units
- Syntax and semantics provide a language's definition
  - Users of a language definition
    - Other language designers
    - Implementers
    - Programmers (the users of the language)



New Jersey's Science & Technology University

COLLEGE OF COMPUTING SCIENCES

## Terminology

- A sentence is a string of characters over some alphabet
- A language is a set of sentences
- A lexeme is the lowest level syntactic unit of a language (e.g., \*, sum, begin)
- A token is a category of lexemes (e.g., identifier)



New Jersey's Science & Technology University

COLLEGE OF COMPUTING SCIENCES

## Formal Definition of Languages

- Recognizers
  - A recognizer reads input strings over the alphabet of the language and decides whether the input strings belong to the language
  - The syntax analysis part of a compiler – the parser - is a recognizer

## BNF and Context-Free Grammars

- Context-Free Grammars
  - Developed by Noam Chomsky in the mid-1950s
  - Define a class of languages called context-free languages
- Backus-Naur Form (1959)
  - Invented by John Backus to describe the syntax of Algol 58
  - BNF is used to represent context-free grammars

## BNF Fundamentals

- In BNF, abstractions are used to represent classes of syntactic structures--they act like syntactic variables (also called nonterminal symbols, or just nonterminals)
- Terminals are lexemes or tokens
- A rule has a left-hand side (LHS), which must be a single nonterminal, and a right-hand side (RHS), which is a string of terminals and/or nonterminals

## BNF Fundamentals (continued)

- Examples of BNF rules:
 

```
<ident_list> → identifier | identifier comma <ident_list>
<if_stmt> → if <logic_expr> then <stmt>

ident_list ::= IDENTIFIER | IDENTIFIER COMMA ident_list
is_stmt ::= IF logic_expr THEN stmt
```
- Grammar: a finite non-empty set of rules
- The start symbol is a special element of the nonterminals of a grammar; it indicates where the recognition of the language begins

## BNF Rules

- A nonterminal symbol can have more than one RHS

```
<stmt> → <single_stmt>  
        | begin <stmt_list> end
```

- Syntactic lists are described using recursion

```
<ident_list> → ident  
              | ident, <ident_list>
```

## Derivation

- A derivation is a repeated application of rules, starting with the start symbol and ending with a sentence (all terminal symbols)

## An Example Grammar

1.  $\langle \text{program} \rangle \rightarrow \langle \text{stmts} \rangle$
2.  $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$
3.  $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
4.  $\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid d$
5.  $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle - \langle \text{term} \rangle$
6.  $\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \text{const}$



New Jersey's Science &amp; Technology University

COLLEGE OF COMPUTING SCIENCES

## An Example Derivation

For this sentence:

$a = b + \text{const}$

1.  $\langle \text{program} \rangle \rightarrow \langle \text{stmts} \rangle$
2.  $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$
3.  $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
4.  $\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid d$
5.  $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle - \langle \text{term} \rangle$
6.  $\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \text{const}$

```

<program>  => <stmts>      # rule 1
            => <stmt>      # rule 2
            => <var> = <expr> # rule 3
            => a = <expr>    # rule 4
            => a = <term> + <term> # rule 5
            => a = <var> + <term> # rule 6
            => a = b + <term> # rule 4
            => a = b + const  # rule 6
  
```



New Jersey's Science &amp; Technology University

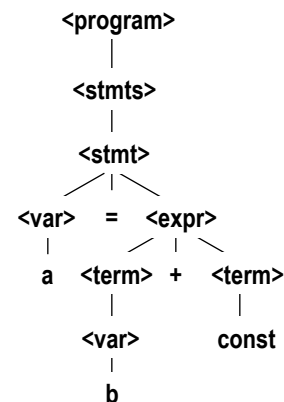
COLLEGE OF COMPUTING SCIENCES

## About Derivations

- Every string of symbols in a derivation is a sentential form
- A sentence is a sentential form that has only terminal symbols
- A leftmost derivation is one in which the leftmost nonterminal in each sentential form is the one that is expanded
- A derivation may be either leftmost or rightmost

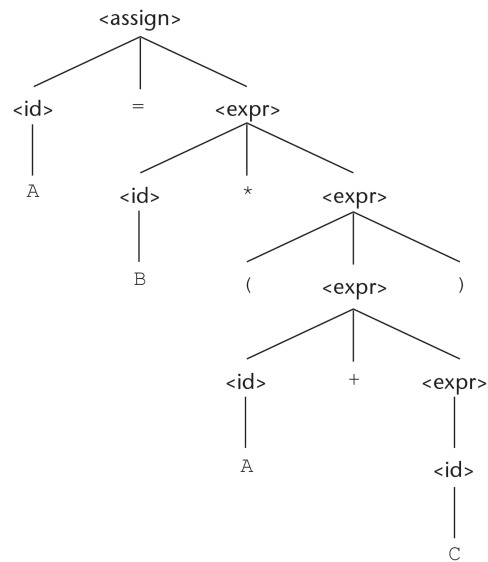
## Parse Tree

- A hierarchical representation of a derivation
- Each step in the derivation is a level in the tree
- A traversal of the tree is a representation of the expression that you parsed



### A parse tree for the simple statement

**A = B \* (A + C)**



## Ambiguity in Grammars

- A grammar is ambiguous if it generates a sentential form that has two or more distinct parse trees

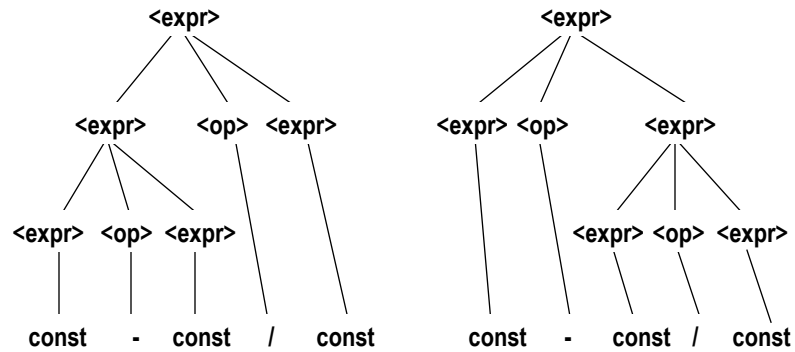


## An Ambiguous Expression Grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{const}$

$\langle \text{op} \rangle \rightarrow / \mid -$

$\text{const} - \text{const} / \text{const}$



**NJIT**

New Jersey's Science & Technology University

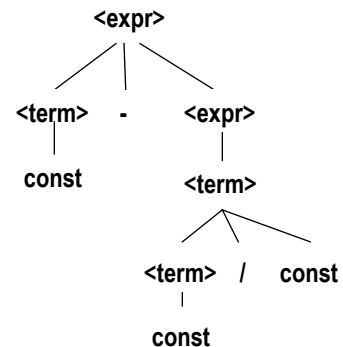
COLLEGE OF COMPUTING SCIENCES

## An Unambiguous Expression Grammar

- If we use the parse tree to indicate precedence levels of the operators, we will not have ambiguity

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle - \langle \text{expr} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle / \text{const} \mid \text{const}$



**NJIT**

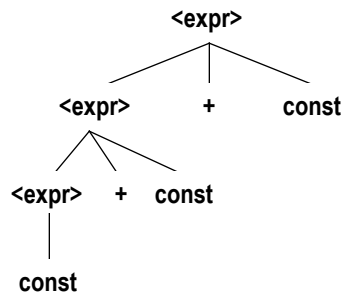
New Jersey's Science & Technology University

COLLEGE OF COMPUTING SCIENCES

## Associativity of Operators

- Operator associativity can also be indicated by a grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \text{const}$  (ambiguous)  
 $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \text{const} \mid \text{const}$  (unambiguous)



## Extended BNF Notation

- Optional parts are placed in brackets [ ]  
 $\langle \text{proc\_call} \rangle \rightarrow \text{ident } ([\langle \text{expr\_list} \rangle])$
- Alternative parts of RHSs are placed inside parentheses and separated via vertical bars  
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle (+|-) \text{const}$
- Repetitions (0 or more) are placed inside braces { }  
 $\langle \text{ident} \rangle \rightarrow \text{letter } \{\text{letter}|\text{digit}\}$

## BNF and EBNF

- BNF

```

<expr> → <expr> + <term>
        | <expr> - <term>
        | <term>
<term> → <term> * <factor>
        | <term> / <factor>
        | <factor>
  
```

- EBNF

```

<expr> → <term> { (+ | -) <term> }
<term> → <factor> { (* | /) <factor> }
  
```

## Variations in EBNF

- Alternative RHSs are put on separate lines
- Use of a colon instead of  $\Rightarrow$
- Use of <sub>opt</sub> for optional parts
- Use of <sub>oneof</sub> for choices

## The Parsing Problem

- Goals of the parser, given an input program:
  - Find all syntax errors; for each, produce an appropriate diagnostic message and recover quickly
  - Produce the parse tree, or at least a trace of the parse tree, for the program

## The Parsing Problem (continued)

- Two categories of parsers
  - Top down - produce the parse tree, beginning at the root
    - Start from the top, work down and to the right
  - Bottom up - produce the parse tree, beginning at the leaves
    - Start from the bottom, work up and to the right
- Useful parsers look only one token ahead in the input

## The Parsing Problem (continued)

- Top-down Parsers
  - Given a sentential form,  $xA\alpha$ , the parser must choose the correct A-rule to get the next sentential form in the leftmost derivation, using only the first token produced by A
- The most common top-down parsing algorithms:
  - Recursive descent - a coded implementation
  - LL parsers - table driven implementation

## The Parsing Problem (continued)

- Bottom-up parsers
  - Given a right sentential form,  $\alpha$ , determine what substring of  $\alpha$  is the right-hand side of the rule in the grammar that must be reduced to produce the previous sentential form in the right derivation
  - The most common bottom-up parsing algorithms are in the LR family

## The Parsing Problem (continued)

- The Complexity of Parsing
  - Parsers that work for any unambiguous grammar are complex and inefficient (  $O(n^3)$ , where  $n$  is the length of the input )
  - Compilers use parsers that only work for a subset of all unambiguous grammars, but do it in linear time (  $O(n)$ , where  $n$  is the length of the input )

## Recursive-Descent Parsing

- There is a subprogram for each nonterminal in the grammar, which can parse sentences that can be generated by that nonterminal
- EBNF is ideally suited for being the basis for a recursive-descent parser, because EBNF minimizes the number of nonterminals

## Recursive-Descent Parsing

- Example: a grammar for simple expressions:

```
<expr> → <term> { (+ | -) <term> }  
<term> → <factor> { (* | /) <factor> }  
<factor> → id | int_constant | ( <expr> )
```

## Recursive-Descent Parsing

- Example assumes we have a lexical analyzer in a function named `lex`, which puts the next token code in `nextToken`
- Example assumes we have a function named `putBackToken`, which pushes back the last token read
- NOTE: most error detection and error handling is NOT shown in this example

## Recursive-Descent Parsing

- The coding process when there is only one RHS in a rule:
  - For each terminal symbol in the RHS, compare it with the next input token; if they match, continue, else there is an error
  - For each nonterminal symbol in the RHS, call its associated parsing subprogram

## Recursive-Descent Parsing

```

/* Function expr
   Parses strings in the language
   generated by the rule:
   <expr> -> <term> { (+ | -) <term> }
*/

void expr() {
    /* Parse the first term */
    term(); // IF THIS FAILS, THEN THIS IS NOT AN expr(): Error

    /* As long as the next token is + or -, call
       * lex to get the next token and parse the
       * next term */
    while (true) {
        lex();
        if (nextToken == ADD_OP || nextToken == SUB_OP) {
            term(); // IF THIS FAILS? error
        }
        else {
            putBackToken();
            break;
        }
    }
}

```



## Recursive-Descent Parsing

```

/* term
  Parses strings in the language generated by the rule:
  <term> -> <factor> { (* | /) <factor> }
*/

void term() {
    /* Parse the first factor */
    factor(); // IF THIS FAILS, THEN THIS IS NOT AN expr(): Error

    /* As long as the next token is * or /, call
    * lex to get the next token and parse the
    * next factor */
    while (true) {
        lex();
        if (nextToken == MULT_OP || nextToken == DIV_OP) {
            factor(); // IF THIS FAILS? error
        }
        else {
            putBackToken();
            break;
        }
    }
}

```

## Recursive-Descent Parsing

- A nonterminal that has more than one RHS requires an initial process to determine which RHS it is to parse
  - The correct RHS is chosen on the basis of the next token of input (the lookahead)
  - The next token is compared with the first token that can be generated by each RHS until a match is found
  - If no match is found, it is a syntax error

# Recursive-Descent Parsing

```

/* Function factor
Parses strings in the language
generated by the rule:
<factor> -> id | int_constant | (<expr>) */

void factor() {
    lex();

    /* Determine which RHS */
    if (nextToken == ID_CODE || nextToken == INT_CODE )
        return; // SUCCESS

    /* See if the RHS is (<expr>) */
    else if (nextToken == LP_CODE) {
        expr(); // IF FAILS.. error
        lex();
        if (nextToken == RP_CODE)
            return; // SUCCESS
        else
            error(); // mismatched paren
    }
    else error(); /* Neither RHS matches */
}

```

# Recursive-Descent Parsing

- Trace of (sum + 47) / total

```

Call expr()
  Call term()
    Call factor()
      Token is LP_CODE:(
      Call expr()
        Call term()
          Call factor()
            Token is ID_CODE:sum
            Return from factor
          Return from term
        Token is ADD_OP:+
        Call term()
          Call factor()
            Token is INT_CODE:47
            Return from factor
          Return from term
        Return from expr
      Token is RP_CODE:)
    Return from factor
  Token is DIV_OP:/
  Call factor()
    Token is ID_CODE:total
  Return from factor
Return from term
Return from expr

```

## Recursive-Descent Parsing

- The LL Grammar Class
  - The Left Recursion Problem
    - If a grammar has left recursion, either direct or indirect, it cannot be the basis for a top-down parser
      - A grammar can be modified to remove direct left recursion as follows:  
For each nonterminal, A,
        1. Group the A-rules as  $A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$   
where none of the  $\beta$ 's begins with A
        2. Replace the original A-rules with  
 $A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$   
 $A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$

## Recursive-Descent Parsing

- The other characteristic of grammars that disallows top-down parsing is the lack of pairwise disjointness
  - The inability to determine the correct RHS on the basis of one token of lookahead
  - Def:  $\text{FIRST}(\alpha) = \{a \mid \alpha \Rightarrow^* a\beta\}$   
(If  $\alpha \Rightarrow^* \epsilon$ ,  $\epsilon$  is in  $\text{FIRST}(\alpha)$ )

## Recursive-Descent Parsing

- Pairwise Disjointness Test:
  - For each nonterminal, A, in the grammar that has more than one RHS, for each pair of rules,  $A \rightarrow \alpha_i$  and  $A \rightarrow \alpha_j$ , it must be true that
 
$$\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \phi$$
- Example:
 
$$A \rightarrow a \mid bB \mid cAb$$

$$A \rightarrow a \mid aB$$

## Recursive-Descent Parsing

- Left factoring can resolve the problem
- Replace
- $$\langle \text{variable} \rangle \rightarrow \text{identifier} \mid \text{identifier} [\langle \text{expression} \rangle]$$
- with
- $$\langle \text{variable} \rangle \rightarrow \text{identifier} \langle \text{new} \rangle$$
- $$\langle \text{new} \rangle \rightarrow \epsilon \mid [\langle \text{expression} \rangle]$$
- or
- $$\langle \text{variable} \rangle \rightarrow \text{identifier} [[\langle \text{expression} \rangle]]$$

## Bottom-up Parsing

- Can be used with left-recursive grammars
- The parsing problem is finding the correct RHS in a right-sentential form to reduce to get the previous right-sentential form in the derivation

## Example

$E ::= E + T \mid T$   
 $T ::= T * F \mid F$   
 $F ::= ( E ) \mid \text{id}$

Rightmost derivation:

$E \Rightarrow E + T$   
 $\Rightarrow E + T * F$   
 $\Rightarrow E + T * \text{id}$   
 $\Rightarrow E + F * \text{id}$   
 $\Rightarrow E + \text{id} * \text{id}$   
 $\Rightarrow T + \text{id} * \text{id}$   
 $\Rightarrow F + \text{id} * \text{id}$   
 $\Rightarrow \text{id} + \text{id} * \text{id}$

## Bottom-up Parsing (continued)

- $E + T * id$ 
  - contains 3 right hand sides:  $E+T$ ,  $T$ , and  $id$
  - The one that should be chosen is the “handle”: the right hand side that should be rewritten to get the next sequential form in the derivation

## Bottom-up Parsing (continued)

- Handles:
  - Def:  $\beta$  is the *handle* of the right sentential form  $\gamma = \alpha\beta w$  if and only if  $S \Rightarrow_{rm}^* \alpha A w \Rightarrow_{rm} \alpha\beta w$
  - Def:  $\beta$  is a *phrase* of the right sentential form  $\gamma$  if and only if  $S \Rightarrow_{rm}^* \gamma = \alpha_1 A \alpha_2 \Rightarrow_{rm} \alpha_1 \beta \alpha_2$
  - Def:  $\beta$  is a *simple phrase* of the right sentential form  $\gamma$  if and only if  $S \Rightarrow_{rm}^* \gamma = \alpha_1 A \alpha_2 \Rightarrow_{rm} \alpha_1 \beta \alpha_2$

## Bottom-up Parsing (continued)

- Intuition about handles:
  - The handle of a right sentential form is its leftmost simple phrase
  - Given a parse tree, it is now easy to find the handle
  - Parsing can be thought of as handle pruning

## Bottom-up Parsing (continued)

- Shift-Reduce Algorithms
  - Reduce is the action of replacing the handle on the top of the parse stack with its corresponding LHS
  - Shift is the action of moving the next token to the top of the parse stack

## Bottom-up Parsing (continued)

- Advantages of LR parsers:
  - They will work for nearly all grammars that describe programming languages.
  - They work on a larger class of grammars than other bottom-up algorithms, but are as efficient as any other bottom-up parser.
  - They can detect syntax errors as soon as it is possible.
  - The LR class of grammars is a superset of the class parsable by LL parsers.

## Bottom-up Parsing (continued)

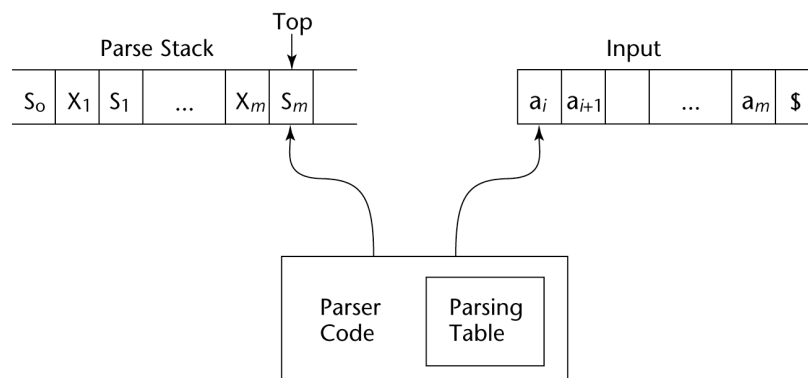
- A parser table can be generated from a given grammar with a tool, e.g., **yacc** or **bison**
- Knuth's insight: A bottom-up parser could use the entire history of the parse, up to the current point, to make parsing decisions
  - There are only a finite and relatively small number of different parse situations that could have occurred, so the history could be stored in a parser state, on the parse stack



## Bottom-up Parsing (continued)

- LR parsers are table driven, where the table has two components, an ACTION table and a GOTO table
  - The ACTION table specifies the action of the parser, given the parser state and the next token
    - Rows are state names; columns are terminals
  - The GOTO table specifies which state to put on top of the parse stack after a reduction action is done
    - Rows are state names; columns are nonterminals

## Structure of An LR Parser



## Bottom-up Parsing (continued)

- Initial configuration:  $(S_0, a_1 \dots a_n \$)$
- Parser actions:
  - For a Shift, the next symbol of input is pushed onto the stack, along with the state symbol that is part of the Shift specification in the Action table
  - For a Reduce, remove the handle from the stack, along with its state symbols. Push the LHS of the rule. Push the state symbol from the GOTO table, using the state symbol just below the new LHS in the stack and the LHS of the new rule as the row and column into the GOTO table

## Bottom-up Parsing (continued)

- Parser actions (continued):
  - For an Accept, the parse is complete and no errors were found.
  - For an Error, the parser calls an error-handling routine.

## Example Grammar

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T * F$
4.  $T \rightarrow F$
5.  $F \rightarrow (E)$
6.  $F \rightarrow id$

**NJIT**

New Jersey's Science & Technology University

COLLEGE OF COMPUTING SCIENCES

## LR Parsing Table

State	Action						Goto		
	id	+	*	(	)	\$	E	T	F
0	S5		S4				1	2	3
1		S6				accept			
2		R2	S7		R2	R2			
3		R4	R4		R4	R4			
4	S5			S4			8	2	3
5		R6	R6		R6	R6			
6	S5			S4				9	3
7	S5			S4					10
8		S6			S11				
9		R1	S7		R1	R1			
10		R3	R3		R3	R3			
11		R5	R5		R5	R5			

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T * F$
4.  $T \rightarrow F$
5.  $F \rightarrow (E)$
6.  $F \rightarrow id$

stack	input	action
0	id+id*id	S5
0id5	+id*id	R6 (0,F)
0F3	+id*id	R4 (0,T)
0T2	+id*id	R2 (0,E)
0E1	+id*id	S6
0E1+6	id*id	S5
0E1+6id5	*id	R6 (6,F)
0E1+6F3	*id	R4 (6,T)
0E1+6T9	*id	S7
0E1+6T9*7	id	S5
0E1+6T9*7id5		R6 (7,F)
0E1+6T9*7F10		R3 (6,T)
0E1+6T9		R1 (0,E)
0E1		accept

**NJIT**

New Jersey's Science & Technology University

COLLEGE OF COMPUTING SCIENCES

## Bottom-up Parsing (continued)

- If, when the table is constructed, two entries must be made in one slot in the table, this is an ambiguity (a “shift-reduce” conflict or a “reduce-reduce” conflict)

## Resources

- Leftmost vs Rightmost Derivation
  - <https://www.seas.upenn.edu/~cit596/notes/dave/cfg8.html>
  - Note: all three ended up with “abbc”, but they didn’t HAVE to.
- Handles
  - <http://sites.tufts.edu/comp181/2013/10/06/shift-reduce-parsing-bottom-up-parsing>
- LR Parsing Tables
  - Skip to the full table on page 3, unless you want to go through the process of building the table:
  - <http://web.eecs.umich.edu/~weimerw/2009-4610/reading/lr-example.pdf>

